

22 Febbraio 2019 - Analisi Esplorativa

Cognome:

Nome:

Matricola:

Tipologia d'esame: 12 CFU 15 CFU

Prova scritta - fila A

Si svolgano gli esercizi riportando il risultato dove indicato. Durata: 80 minuti

Esercizio 1 (10 punti)

Il dataset `Cars93` presente nella libreria `MASS` descrive 93 modelli di automobile rispetto a 27 variabili. Le variabili di interesse per l'analisi sono $p = 12$ come elencato nel seguito:

- `Weight` peso (in libbre)
- `Width` larghezza (in pollici)
- `Length` lunghezza (in pollici)
- `EngineSize` cilindrata (in litri)
- `Horsepower` potenza
- `RPM` giri al minuto
- `Fuel.tank.capacity` capacità di carburante
- `Passengers` numero di passeggeri
- `MPG.highway` Miglia per gallone in autostrada
- `MPG.city` Miglia per gallone in città
- `Price` prezzo (in migliaia di dollari)
- `Wheelbase` interasse (in pollici)

Sia X la matrice 93×12 corrispondente al dataset `Cars93` ridotto selezionando le variabili sopra elencate.

- a. Si calcoli $d_M^2(x_i, \bar{x})$, il quadrato della distanza di Mahalanobis di ciascuna auto (ciascuna riga della matrice X) dal baricentro. Si riportino i valori di $d_M^2(x_i, \bar{x})$ solo se superano il valore 21, specificando anche il nome dell'automobile a cui si fa riferimento (informazione reperibile dalla variabile `Model` del dataset `Cars93`).

```
rm(list=ls())
library(MASS)
vars = c("Weight", "Width", "Length", "EngineSize", "Horsepower", "RPM", "Fuel.tank.capacity", "Passengers")
X = Cars93[,vars]
n = nrow(X)
p = ncol(X)
row.names(X) <- Cars93[, "Model"]
xbar = matrix(colMeans(X), nrow=p, ncol=1)
S = var(X) * ((n-1)/n)
```

```
InvS = solve(S)
dM2 = apply(X,MARGIN=1, function(u) t(u-xbar) %*% InvS %*% (u - xbar) )
round(dM2[dM2 > 21],2)
```

Astro	Corvette	Stealth	Aerostar	Metro	Civic	RX-7	300E
34.26	39.02	38.86	26.44	28.58	28.02	42.64	35.97

- b. Svolgere l'analisi delle componenti principali sui dati standardizzati Z che si ottengono a partire da X . Si riportino i punteggi (*scores*) per l'automobile Firebird relativamente alle componenti principali 2, 8 e 10 (arrotondati alla terza cifra decimale).

```
Z = scale(X,center=T,scale=diag(S)^(1/2))
round(prcomp(Z)$x["Firebird",c(2,8,10)],3)
```

PC2	PC8	PC10
0.128	-0.194	-0.029

- c. Applicare l'algoritmo delle K medie (`algorithm = "Hartigan-Wong"`) per i dati standardizzati Z iniziando i K centri utilizzando le prime K osservazioni (righe $1, \dots, K$ della matrice Z). Arrotondando il risultato alla seconda cifra decimale, riportare per $K = 2, 3, 4, 5$
- il valore dell'indice $CH(K) = \frac{B/(K-1)}{W/(n-K)}$ di Calinski and Harabasz
 - il valore medio della *silhouette* considerando come matrice delle distanze quella ottenuta con la metrica di Lagrange basata su Z

K	2	3	4	5
$CH(K)$				
<i>silhouette</i> (K)				

```
D = dist(Z, method="maximum")
K = 2:5
CH <- vector()
sil <- vector()
library(cluster)
for (k in 1:length(K)){
  km = kmeans(Z, centers=Z[1:K[k],])
  CH[k] = (km$betweenss/(K[k]-1))/(km$tot.withinss/(n-K[k]))
  sil[k] = summary(silhouette(x=km$cluster, dist=D))$avg.width
}
round(rbind(K,CH,sil),2)
```

	[,1]	[,2]	[,3]	[,4]
K	2.00	3.00	4.00	5.00
CH	75.77	61.10	57.01	48.59
sil	0.26	0.18	0.22	0.20

- d. Si consideri la stima di massima verosimiglianza per il modello fattoriale con k fattori basato sui dati standardizzati Z . Riportare il p -value del primo test non significativo al livello 5% (e il corrispondente valore di k) per la sequenza di ipotesi nulle $H_0(k = 1), H_0(k = 2), H_0(k = 3), \dots$ dove $H_0(k)$ ="il modello fattoriale con k fattori è corretto".

```
k = which.max(sapply(1:6, function(k) factanal(Z,factors=k)$PVAL > 0.05 ))
k
```

```
objective
      6
```

```
round(factanal(Z,factors=k)$PVAL,4)
```

```
objective
0.1693
```

e. Stimare il modello fattoriale con $k = 2$ fattori con il metodo della massima verosimiglianza utilizzando i dati standardizzati Z e senza effettuare alcuna rotazione. Riportare il valore della statistica test rapporto di verosimiglianza $T = n \log \left(\frac{|\hat{\Lambda}' + \hat{\Psi}|}{|R|} \right)$ (arrotondando al terzo decimale)

```
R = cor(X)
af = factanal(Z,factors=2, rotation="none", method="mle")
Lambda = af$loadings[,]
Psi = diag(af$uniqueness)
fit = Lambda %*% t(Lambda) + Psi
lrt = n*log(det(fit)/det(R))
round(lrt,3)
```

```
[1] 345.353
```

Esercizio 2 (punti 2)

Dimostrare, esplicitando tutti i passaggi, e specificando tutte le quantità coinvolte.

Se esiste un $c_{p \times 1} \neq 0$ tale che $S_{p \times p} c_{p \times 1} = 0_{p \times 1}$, allora le colonne di $\tilde{X}_{n \times p}$ sono linearmente dipendenti.

Esercizio 3 (3 punti)

Si consideri la seguente matrice di correlazione $R_{3 \times 3} = \begin{bmatrix} 1 & 1/2 & 1/2 \\ 1/2 & 1 & 2/3 \\ 1/2 & 2/3 & 1 \end{bmatrix}$.

- a. Sia \tilde{X} la matrice dei dati centrati. Determinare l'angolo (espresso in gradi) tra \tilde{x}_1 e \tilde{x}_3 : =
- b. Sapendo che $s_{11} = 4$, $s_{22} = 9$ e $\text{tr}(S) = 14$, calcolare la matrice di varianze/covarianze S e $S^{1/2}$ (arrotondando al secondo decimale)

$$S = \begin{bmatrix} \dots & \dots & \dots \\ \dots & \dots & \dots \\ \dots & \dots & \dots \end{bmatrix} \quad S^{1/2} = \begin{bmatrix} \dots & \dots & \dots \\ \dots & \dots & \dots \\ \dots & \dots & \dots \end{bmatrix}.$$

```
R = matrix(c(1,1/2,1/2,1/2,1,2/3,1/2,2/3,1), byrow=T, ncol=3)
```

```
# a
acos(R[1,3])*(180/pi)
```

```
[1] 60
```

```
# b
s11=4
s22=9
s33=14-s11-s22
D = diag(c(s11,s22,s33))
S = D^(1/2) %*% R %*% D^(1/2)
S
```

```
      [,1] [,2] [,3]
[1,]    4    3    1
[2,]    3    9    2
[3,]    1    2    1
```

```
lambdas = eigen(S)$values
V = eigen(S)$vectors
sqrtS = V %*% diag( lambdas^(1/2) ) %*% t(V)
round(sqrtS,2)
```

```
      [,1] [,2] [,3]
[1,] 1.89  0.6 0.26
[2,] 0.60  2.9 0.50
[3,] 0.26  0.5 0.83
```

Esercizio 4 (2 punti)

Si consideri la seguente matrice di correlazione calcolata sulla base di $n = 50$ osservazioni:

	a	b	c	d
a	1	0.8	0.1	0.6
b	0.8	1	0.3	0.7
c	0.1	0.3	1	0.4
d	0.6	0.7	0.4	1

Sulla base della matrice di correlazione, si stimi il modello fattoriale con $k = 1$ fattori utilizzando il metodo della massima verosimiglianza senza effettuare alcuna rotazione. Si determini il punteggio fattoriale con il metodo di Thompson (arrotondando alla quarta cifra decimale) per una certa unità statistica sapendo che i suoi valori standardizzati nelle quattro variabili a , b , c e d sono

	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>
unità	1.2426	0.7828	-0.5209	-0.0034

```
rm(list=ls())
X = USArrests
names(X) = c("a", "b", "c", "d")
R = round(cor(X), 1)
hatLambda = matrix(factanal(covmat=R, factors=1, n.obs=50, rotation="none")$loadings[, ], ncol=1)
z = c(1.2426, 0.7828, -0.5209, -0.0034)
Rinv = solve(R)
round( t(hatLambda) %*% Rinv %*% z, 4)
```

```
[,1]
[1,] 0.7822
```

$\hat{f} =$

Esercizio 5 (3 punti)

Alla matrice $S_{p \times p}$ sono associati i seguenti autovalori $\lambda_1 = 6, \lambda_2 = 4$ e autovettori normalizzati $v_1 = \begin{bmatrix} 1/\sqrt{5} \\ 2/\sqrt{5} \end{bmatrix}$,

$$v_2 = \begin{bmatrix} 2/\sqrt{5} \\ -1/\sqrt{5} \end{bmatrix}.$$

a. Calcolare la correlazione tra la prima colonna \tilde{x}_1 di \tilde{X} e i punteggi y_1 della prima componente principale, arrotondando al secondo decimale

=

b. Determinare (arrotondando al secondo decimale) la matrice di correlazione

$$R = \begin{bmatrix} & \\ & \end{bmatrix}$$

```
rm(list=ls())
# a.
Lambda = diag(c(6,4))
V = matrix(c(1/sqrt(5),2/sqrt(5),2/sqrt(5),-1/sqrt(5)), byrow=F, ncol=2)
S = V %*% Lambda %*% t(V)
round( V[1,1]*sqrt(Lambda[1,1])/sqrt(S[1,1]) , 2)
```

```
[1] 0.52
```

```
# b.
R = diag(diag(S)^(-1/2)) %*% S %*% diag(diag(S)^(-1/2))
round(R,2)
```

```
[,1] [,2]
[1,] 1.00 0.16
[2,] 0.16 1.00
```

Esercizio 6 (6 punti)

Dimostrare, esplicitando tutti i passaggi, e specificando tutte le quantità coinvolte.

- a. Nel modello fattoriale a k fattori, $\mathbb{E} \begin{pmatrix} x & f' \end{pmatrix} = \Lambda$.
- b. Una generica matrice di varianze/covarianze S è semidefinita positiva.
- c. $\det(S^Y) = \det(S)$ dove $Y = \tilde{X}V$ e le colonne di V sono gli autovettori normalizzati di S